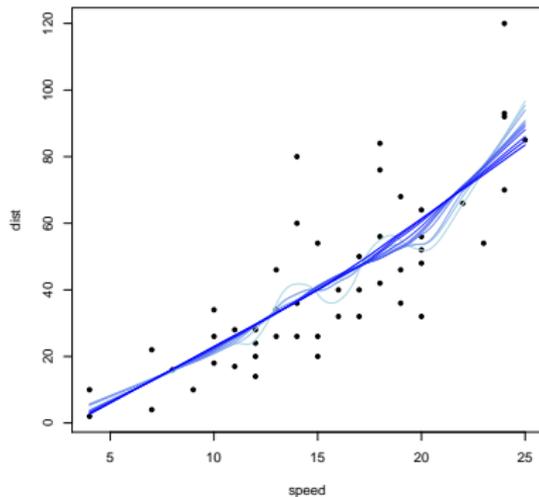


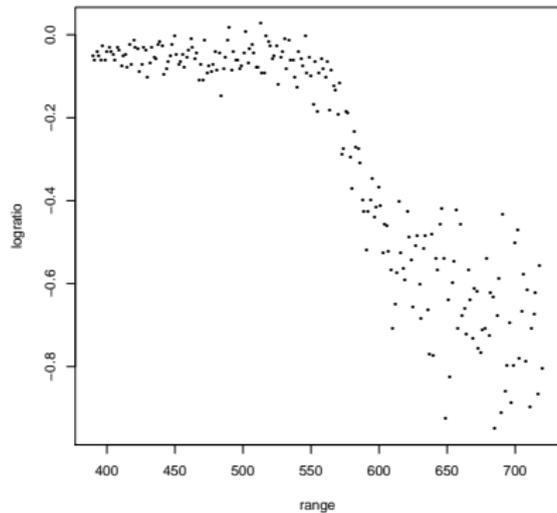
# Modélisation Non Linéaire: Introduction



Yannig Goude

EDF R&D -dpt OSIRIS- Clamart

## *Exemple: lidar data*



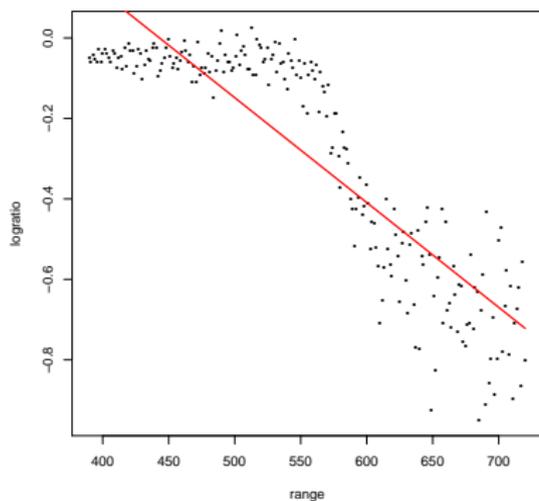
Wikipédia: *light detection and ranging*, est une technologie de télédétection ou de mesure optique basée sur l'analyse des propriétés d'une lumière laser renvoyée vers son émetteur.

Technique utilisée pour détecter des composants chimiques dans l'atmosphère (polluant notamment).

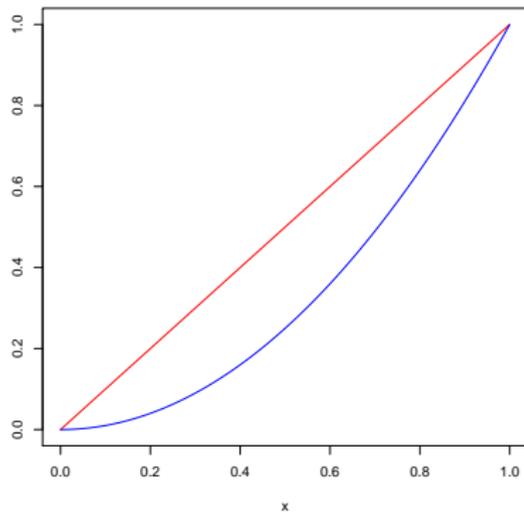
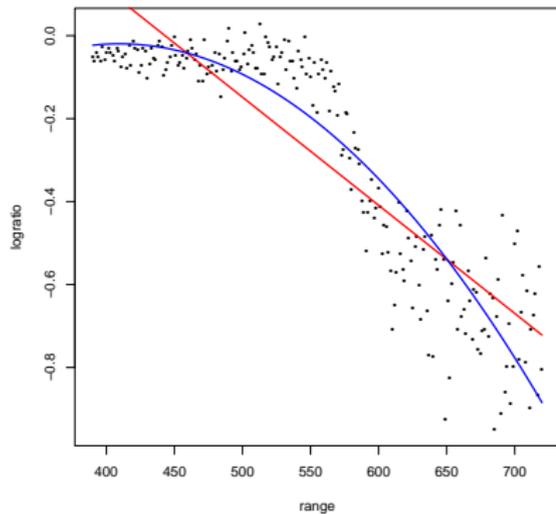
- ▶ la variable *range* est la distance parcourue avant que la lumière revienne au laser
- ▶ la variable *logratio* est le log du rapport de lumière reçue par 2 sources laser. Une source à une fréquence de résonance correspondant au composant d'intérêt -ici le mercure-, l'autre non.

la concentration d'un composant à une distance donnée est donc proportionnelle à la pente de la courbe estimée.

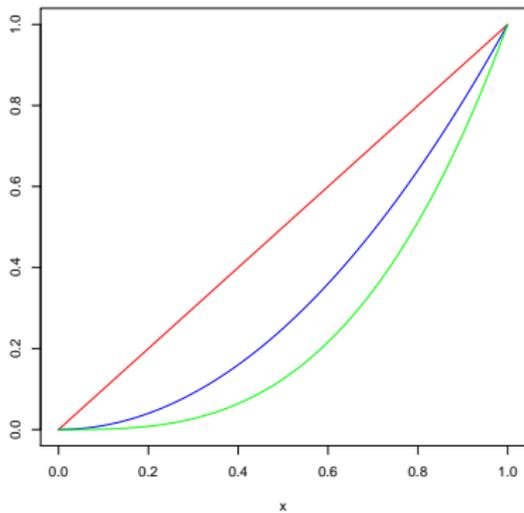
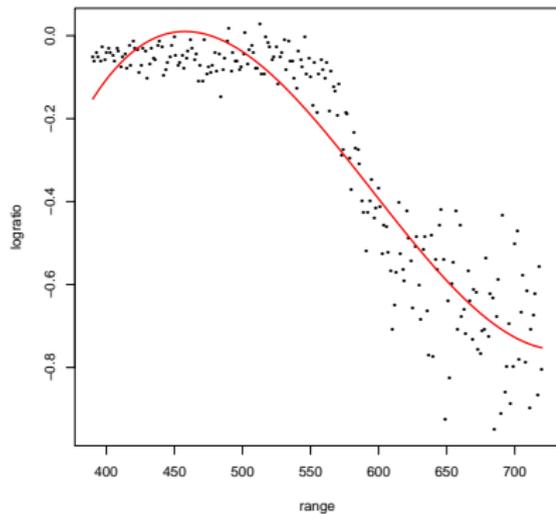
Exemple1: lidar data set, régression linéaire sur  $1, x$



Exemple1: modélisation linéaire sur:  $\mathbf{1}$ ,  $x$  et  $x^2$



Exemple1: trans. polynômiale, reg. linéaire sur:  $\mathbf{1}$ ,  $x$ ,  $x^2$  et  $x^3$



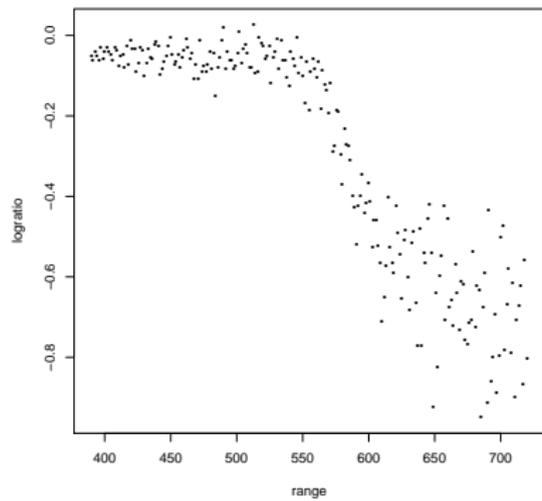
## Kernel Regression

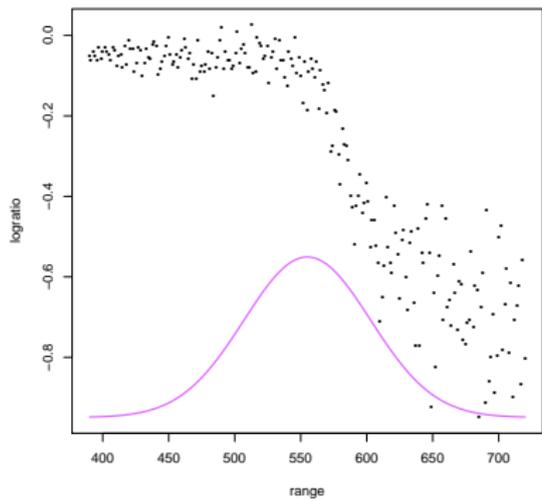
$$\hat{s}(x) = \sum_{i=1}^n w_i(x) y_i$$

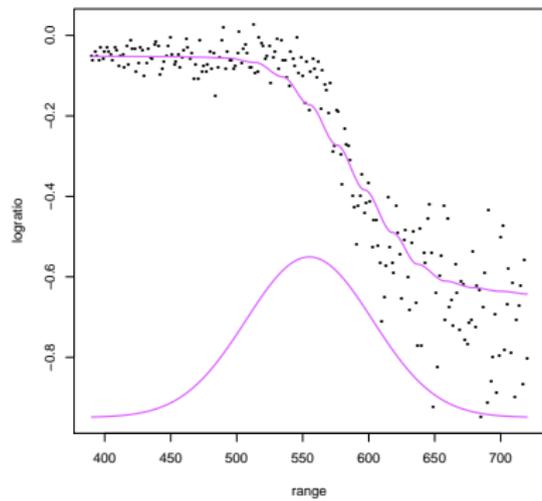
Avec:

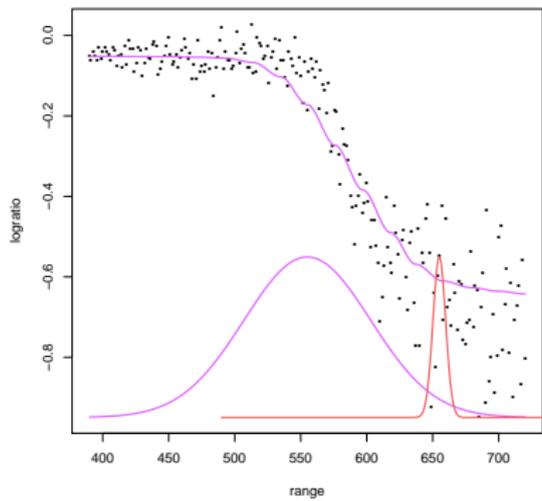
$$w_i(x) \sim \exp(-(x - x_i)^2/h)$$

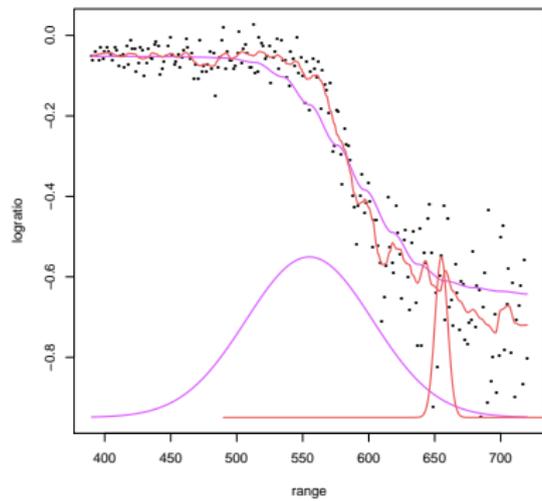
- ▶  $\hat{s}(x)$  est une moyenne pondérée des  $y$
- ▶ plus une observation est proche de  $x$  plus son poids est important
- ▶  $h$ : la "fenêtre" détermine ou le degré de lissage de l'estimateur, sa *régularité*

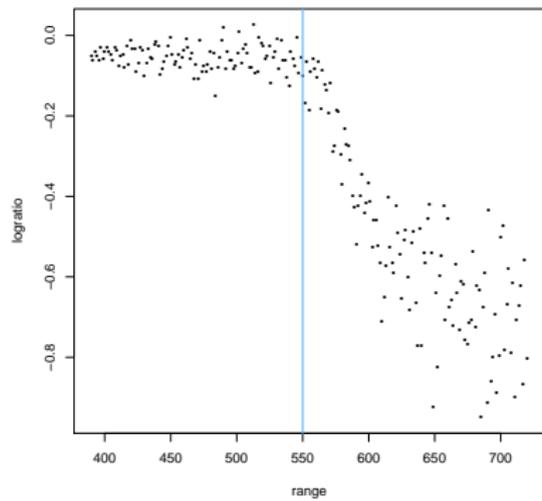


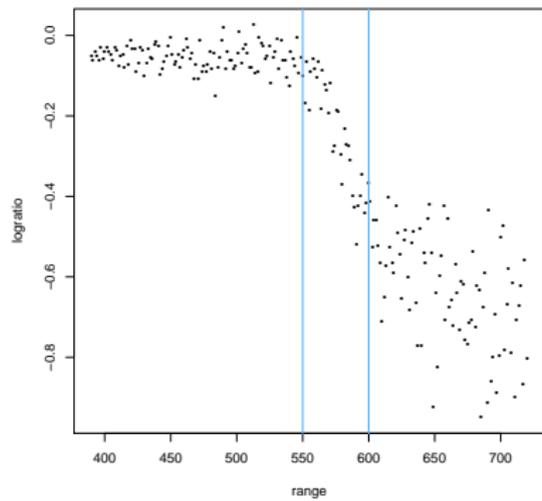




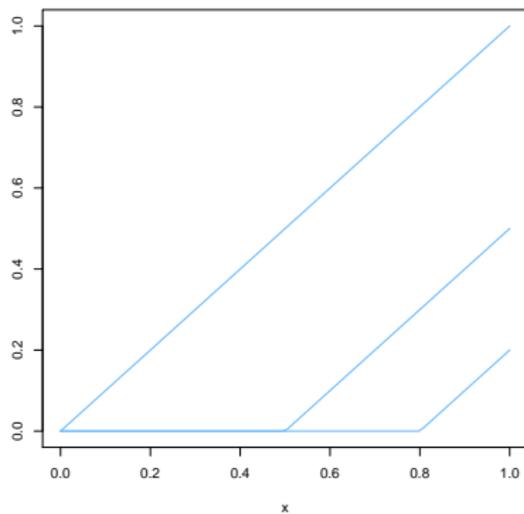
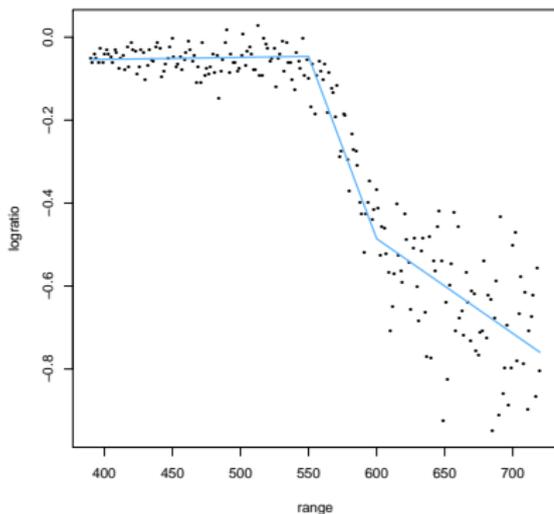






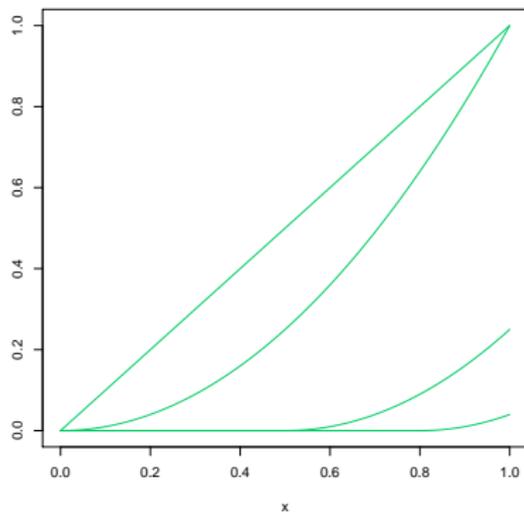
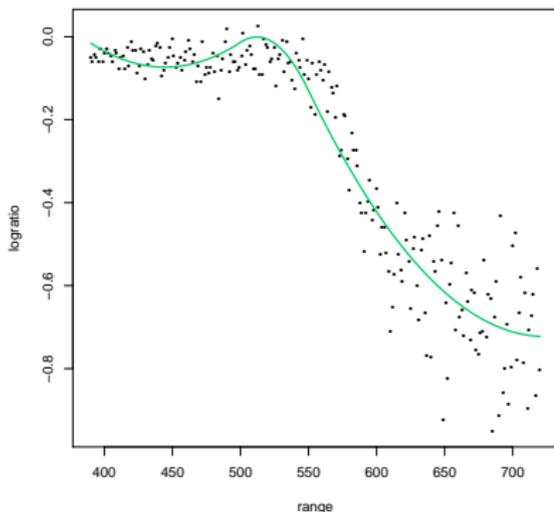


Exemple1: régression linéaire sur  $\mathbf{1}, x, (x - 550)_+, (x - 600)_+$



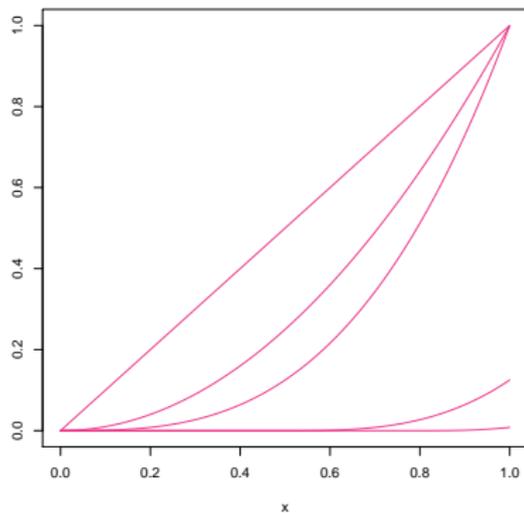
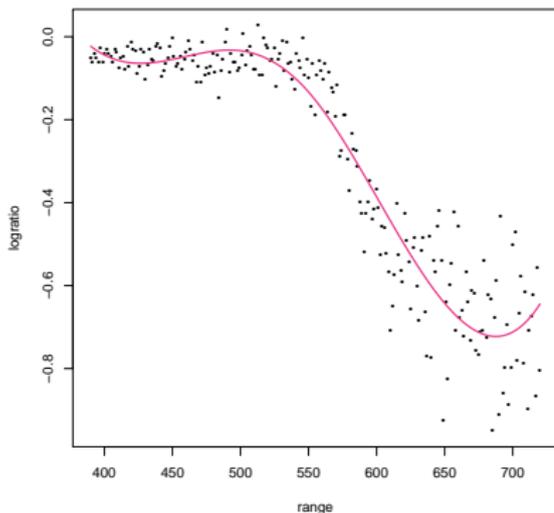
550 et 600 sont appelés des *noeuds*, et les fonctions  $(x - k)_+$  sont appelées *splines linéaires*. L'ensemble des fonctions  $(x - k)_+$  est une base de spline linéaire.

Exemple1: régression linéaire sur  $\mathbf{1}, x, x^2, (x - 550)_+^2, (x - 600)_+^2$



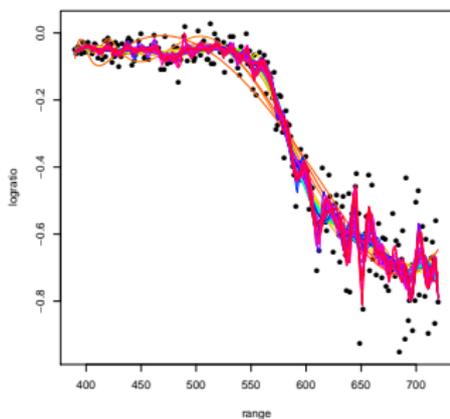
les fonctions  $(x - k)_+^2$  sont appelées *splines quadratiques*

Exemple1: régression linéaire sur  $\mathbf{1}, x, x^2, x^3, (x - 550)_+^3, (x - 600)_+^3$



les fonctions  $(x - k)_+^3$  sont appelées *splines cubiques*

Pb: choix du nombre et de la position des "noeuds", du degré des splines



2 alternatives:

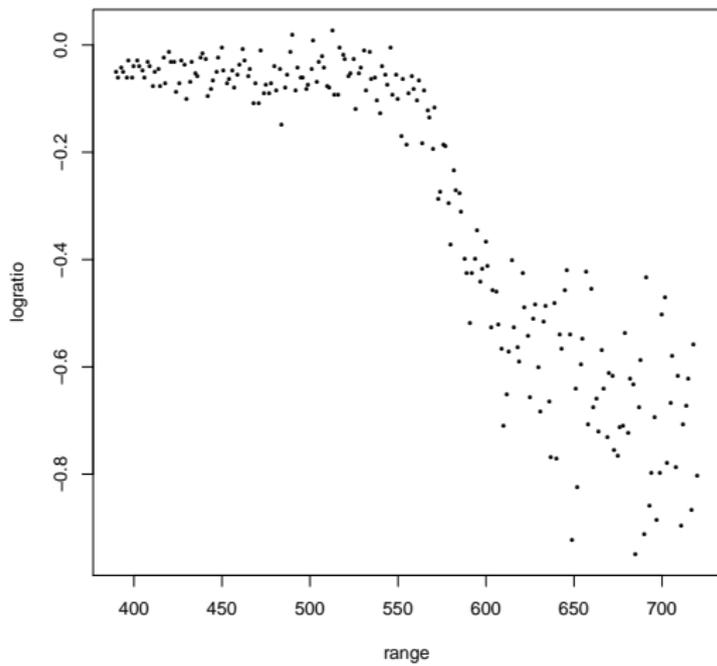
- ▶ Sélection de modèle
- ▶ Méthode de pénalisation

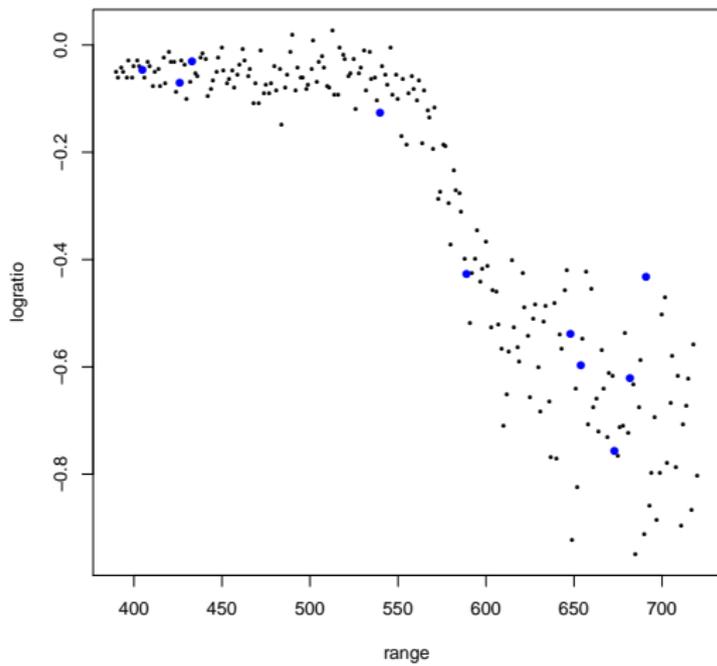
## *Sélection de modèle*

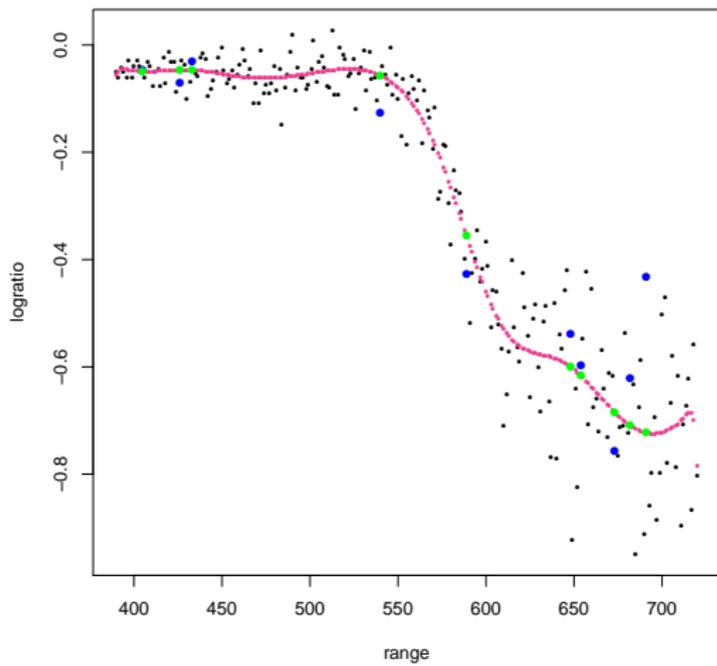
- ▶ Sélection de modèle en régression
  - ▶ Stepwise, forward, backward
- ▶ Échantillon test

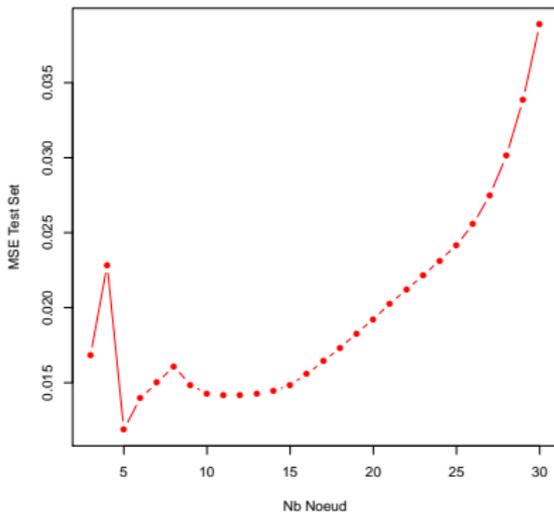
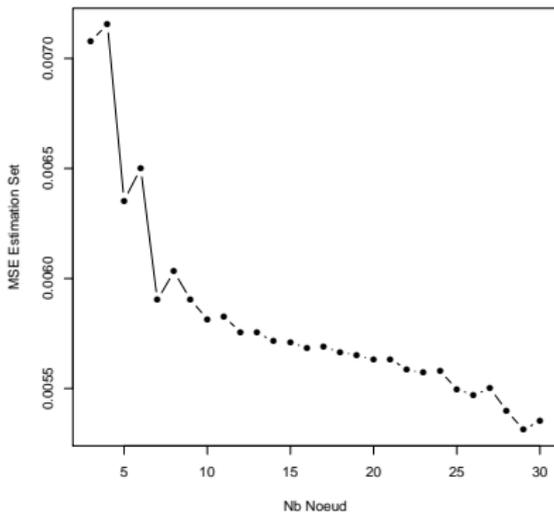
## *Méthode de Pénalisation*

- ▶ AIC, BIC,  $C_p$
- ▶ Pénalisation  $L^1$ : LASSO
- ▶ Pénalisation  $L^2$ : **Ridge regression, P-spline**









⇒ Pour éviter le **surapprentissage**, il faut ajouter une pénalité ou estimer l'erreur de prévision par validation croisée, échantillon test...

